

---

## Learning Probabilities and Causality

### Graded Lab on Variational Autoencoders

---

The lab report is due on December 19th, 2023 and consists of two items:

- The Jupyter Notebook file with your code (to answer section 2.3).
- A PDF with the plots and comments answering sections 2.4 and 2.5.  
This PDF should NOT exceed 3 pages.

Send a ZIP file by e-mail to `xavier.alameda-pineda@inria.fr`.

## 2 Variational Autoencoders for Speech Modeling

### 2.1 Introduction

In this lab work we are going to model the speech signal using variational autoencoders, and derived models. We have prepared the data for you, and you will have access to sequences of magnitude spectrograms. In practice this means that your data will be sequences of  $T = 150$  vectors of dimension  $F = 257$ . These vectors are obtained applying the short-time Fourier transform (which leads to a complex-valued vector) and then computing the entry-wise modulus, therefore we write:  $\mathbf{x}_n \in \mathbb{R}^F$ , where  $n = 1, \dots, N$  is the sample index, and  $N$  is the total number of samples in the training or evaluation sets. Although we have used a large-scale training set to provide a pre-trained model for some of the questions, you will be using two small datasets from two different speakers, Speaker A and Speaker B.

**The generative model** consists of a latent  $Z$ -dimensional variable (also called code), denoted by  $\mathbf{z} \in \mathbb{R}^Z$ . In our experiments we will use  $Z = 16$ . The probabilistic model writes:

$$\mathbf{z} \sim \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I}) \quad \text{and} \quad \mathbf{x}|\mathbf{z} \sim \mathcal{N}_c(\mathbf{x}; \mathbf{0}, \boldsymbol{\nu}_\theta(\mathbf{z})), \quad (1)$$

where  $\mathbf{0}$  and  $\mathbf{I}$  are the zero-valued vector and identity matrix of appropriate dimensions respectively and  $\mathcal{N}_c$  denotes the circular complex Gaussian defined as:

$$p_\theta(\mathbf{x}|\mathbf{z}) = \mathcal{N}_c(\mathbf{x}; \mathbf{0}, \boldsymbol{\nu}_\theta(\mathbf{z})) = \frac{1}{\pi^F |\boldsymbol{\nu}_\theta(\mathbf{z})|} \exp\left(-\mathbf{x}^\top \boldsymbol{\nu}_\theta(\mathbf{z})^{-1} \mathbf{x}\right). \quad (2)$$

where  $\boldsymbol{\theta}$  denote the parameters of the decoder network.

Similarly to the standard VAE formulation seen in class, the covariance matrix depends on the latent code  $\mathbf{z}$ . Differently the standard formulation, the mean is assumed to be zero. In addition, as done in class, we will be implementing a diagonal covariance matrix with log-variances. Formally:

$$\boldsymbol{\nu}_{\boldsymbol{\theta}}(\mathbf{z}) = \begin{pmatrix} \exp(\eta_{1,\boldsymbol{\theta}}(\mathbf{z})) & 0 & \dots \\ 0 & \exp(\eta_{2,\boldsymbol{\theta}}(\mathbf{z})) & \dots \\ \vdots & \vdots & \ddots \end{pmatrix} \quad (3)$$

In other works, the decoder which is the part of the network implementing the generative model, will input  $\mathbf{z}$  and output  $\boldsymbol{\eta}_{\boldsymbol{\theta}}(\mathbf{z}) = [\eta_{1,\boldsymbol{\theta}}(\mathbf{z}), \dots, \eta_{F,\boldsymbol{\theta}}(\mathbf{z})] \in \mathbb{R}^F$ . Recall that  $\boldsymbol{\theta}$  are the parameters of the decoder network.

**The inference model** approximates the posterior of the latent code with a Gaussian distribution:

$$p(\mathbf{z}|\mathbf{x}) \approx q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{\boldsymbol{\phi}}(\mathbf{x}), \boldsymbol{\nu}_{\boldsymbol{\phi}}(\mathbf{x})), \quad (4)$$

where  $\boldsymbol{\phi}$  are the parameters of the so-called encoder network. The covariance matrix of the approximate posterior distribution follows the same structure and rationale as the one of the generative model (i.e. we estimate a vector of log-variances that form a diagonal matrix). As a consequence, the encoder inputs  $\mathbf{x}$  and outputs  $[\boldsymbol{\mu}_{\boldsymbol{\phi}}(\mathbf{x}), \boldsymbol{\eta}_{\boldsymbol{\phi}}(\mathbf{x})] = [\mu_{1,\boldsymbol{\phi}}(\mathbf{x}), \dots, \mu_{Z,\boldsymbol{\phi}}(\mathbf{x}), \eta_{1,\boldsymbol{\phi}}(\mathbf{x}), \dots, \eta_{Z,\boldsymbol{\phi}}(\mathbf{x})] \in \mathbb{R}^{2Z}$ .

**The loss** consists of two terms: reconstruction and regularisation. More formally:

$$\mathcal{L}(\mathbf{x}_{1:N}, \boldsymbol{\theta}, \boldsymbol{\phi}) = \sum_{n=1}^N \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_n|\mathbf{x}_n)} \{ \log p_{\boldsymbol{\theta}}(\mathbf{x}_n|\mathbf{z}_n) \} - \mathcal{KL}(q_{\boldsymbol{\phi}}(\mathbf{z}_n|\mathbf{x}_n) \| p(\mathbf{z}_n)). \quad (5)$$

This is a lower-bound of the log-likelihood that needs to be **maximised**. If you use it within a minimisation framework, you need to change the sign of the expression. Moreover, since the first expectation cannot be taken analytically, we need to approximate with a sample:

$$\mathcal{L}(\mathbf{x}_{1:N}, \boldsymbol{\theta}, \boldsymbol{\phi}) = \sum_{n=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{x}_n|\hat{\mathbf{z}}_n) - \mathcal{KL}(q_{\boldsymbol{\phi}}(\mathbf{z}_n|\mathbf{x}_n) \| p(\mathbf{z}_n)), \quad \hat{\mathbf{z}}_n \sim q_{\boldsymbol{\phi}}(\mathbf{z}_n|\mathbf{x}_n). \quad (6)$$

## 2.2 Pre-implementation work

1. Replace the expression of  $\boldsymbol{\nu}_{\boldsymbol{\theta}}(\mathbf{z})$  into the generative model  $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})$  to compute the sampled log-likelihood  $\log p_{\boldsymbol{\theta}}(\mathbf{x}_n|\hat{\mathbf{z}}_n)$  for the reconstruction term of the loss.
2. Find the formula for the Kullback-Leibler divergence between two multivariate Gaussians, and use it to compute the regularisation term of the loss

$$\mathcal{KL}(q_{\boldsymbol{\phi}}(\mathbf{z}_n|\mathbf{x}_n) \| p(\mathbf{z}_n)), \quad (7)$$

taking into account the structure of  $\boldsymbol{\nu}_{\boldsymbol{\phi}}(\mathbf{x})$ .

3. Write the re-parametrisation trick as a function of the parameters provided by the encoder, that is the mean and the log-variance.

## 2.3 Implementing VAEs

We will be implementing a VAE with the following specifications. The encoder has four dense layers of output dimensions (256, 128, 64, 32) activated with a hyperbolic tangent. After these four layers there is a linear layer to output the mean and the log-variance of  $q_\phi(\mathbf{z}|\mathbf{x})$ . The decoder has also four dense layers with output dimensions (32, 64, 128, 256) activated with a hyperbolic tangent. After these four layers, there is a linear layer to output the log-variance of  $p_\theta(\mathbf{x}|\mathbf{z})$ .

Your implementation work has several steps:

1. Implement the structure of the encoder and decoder in the `build` function. In the code provided, we split the non-linear layers (that should be implemented in `mlp_x_z` and `mlp_z_x`) from the rest.
2. Implement the `generation_x` function (i.e. computing the log-variance of  $p(\mathbf{s}|\mathbf{z})$ ).
3. Implement the `reparametrization` trick as was discussed in class, and you rewrote in the pre-implementation work, in the `reparametrization` function.
4. Implement the `inference` function, to compute the mean and log-variance of  $p(\mathbf{z}|\mathbf{s})$  and sample from it using the `reparametrization` function.
5. Implement the loss with the formulae that you have derived in the pre-implementation work, in the `get_loss` function.

Check your implementation with the short test code provided.

## 2.4 Training and evaluation

In this section we will use the data from Speaker A and Speaker B jointly.

1. Train the implemented architecture for 30 epochs using the provided code. Plot the joint loss as well as the reconstruction and regularisation terms separately.
2. Repeat the training from scratch five times. Plot the average and standard deviation of the loss over the epochs.

## 2.5 Fine-tuning a pre-trained model

In this section we will use the data from Speaker A and Speaker B separately. Download the pre-trained model. For your information, this model has been trained with data collected from around 100 different speakers (which is roughly 9 hours of speech) and for 164 epochs. The pre-trained model is quite different from what you have trained in the previous section.

1. Read the code provided for pre-processing and comment it.
2. Load the model and compute the loss (without any fine-tuning) on the samples of Speaker A and of Speaker B.
3. Fine-tune the pre-trained model on the data of Speaker A and of Speaker B separately, for 20 epochs. You must obtain two models (Model A and Model B). Evaluate each model with the data of the corresponding speaker. What do you observe if you compare it to the previous exercise?
4. Evaluate now Model A with the data of Speaker B, and vice-versa. What do you observe when you compare it with the previous two situations?